# Introduction of VCDE Interoperability Review for Protein Information Resource (PIR)

## caBIG™ VCDE  Workspace Monthly Teleconference
## November 3, 2005

Baris Ethem Suzek
Georgetown University - Lombardi Cancer Center – PIR
bes23@georgetown.edu

# Outline

▶ Introduction

▶ Overview of Grid-Enablement of PIR

▶ Data Model

▶ Semantic Annotation

▶ Using PIR Grid Service

▶ Discussion

# Introduction

▸ **Protein Information Resource (PIR):** Integrated Protein Informatics Resource for Genomic/Proteomic Research



▸ **UniProt Universal Protein Resource:** Central Resource of Protein Sequence and Function

▸ **PIRSF Family Classification System:** Protein Classification and Functional Annotation

▸ **iProClass Integrated Protein Knowledgebase:** Data Integration and Functional Analysis

http://pir.georgetown.edu

# Introduction



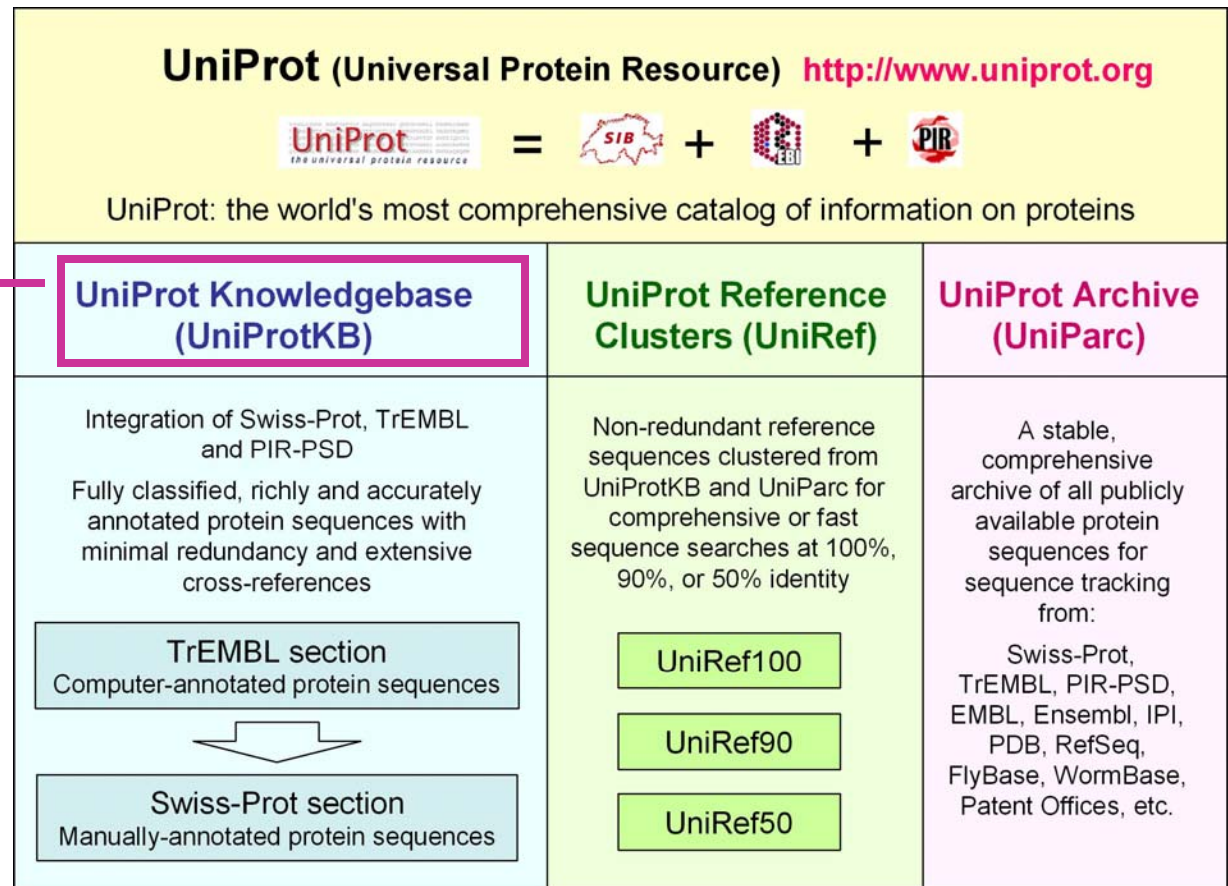▸ UniProt: Universal Protein Resource - Central Resource of Protein Sequence and Function



▸ International Consortium
  – PIR at GUMC
  – European Bioinformatics Institute (EBI)
  – Swiss Institute of Bioinformatics (SIB)

▸ Unifies PIR-PSD, Swiss-Prot, TrEMBL Protein Sequence Databases

http://www.uniprot.org

▶ **UniProt Databases**

Primary data source for Grid-Enablement of PIR



**UniProt** (Universal Protein Resource) http://www.uniprot.org

UniProt = SIB + EBI + PIR

UniProt: the world's most comprehensive catalog of information on proteins

| **UniProt Knowledgebase (UniProtKB)** | **UniProt Reference Clusters (UniRef)** | **UniProt Archive (UniParc)** |
|---|---|---|
| Integration of Swiss-Prot, TrEMBL and PIR-PSD | Non-redundant reference sequences clustered from UniProtKB and UniParc for comprehensive or fast sequence searches at 100%, 90%, or 50% identity | A stable, comprehensive archive of all publicly available protein sequences for sequence tracking from: |
| Fully classified, richly and accurately annotated protein sequences with minimal redundancy and extensive cross-references | | |
| **TrEMBL section** Computer-annotated protein sequences ⬇ **Swiss-Prot section** Manually-annotated protein sequences | UniRef100 UniRef90 UniRef50 | Swiss-Prot, TrEMBL, PIR-PSD, EMBL, Ensembl, IPI, PDB, RefSeq, FlyBase, WormBase, Patent Offices, etc. |

# Project Overview

▶ Goal: Providing methods to query and retrieve protein related information for the cancer research community

▶ Grid-Enablement of PIR project is a data service

▶ All the objects in our model exposed to caGRID as of August 1st

▶ API is generated using caCORE SDK 1.0.3 like caBIO

▶ Example queries:
  – Find the proteins for the gene "BRCA2" (Breast Cancer Gene 2)
  – Find all the proteins that contain the domain BRCA2 repeat (PFAM:PF00634, a domain in Breast cancer type 2 susceptibility protein)
  – ID mapping: Find all the database cross-references from various databases corresponding to RefSeq Accession NP_009225
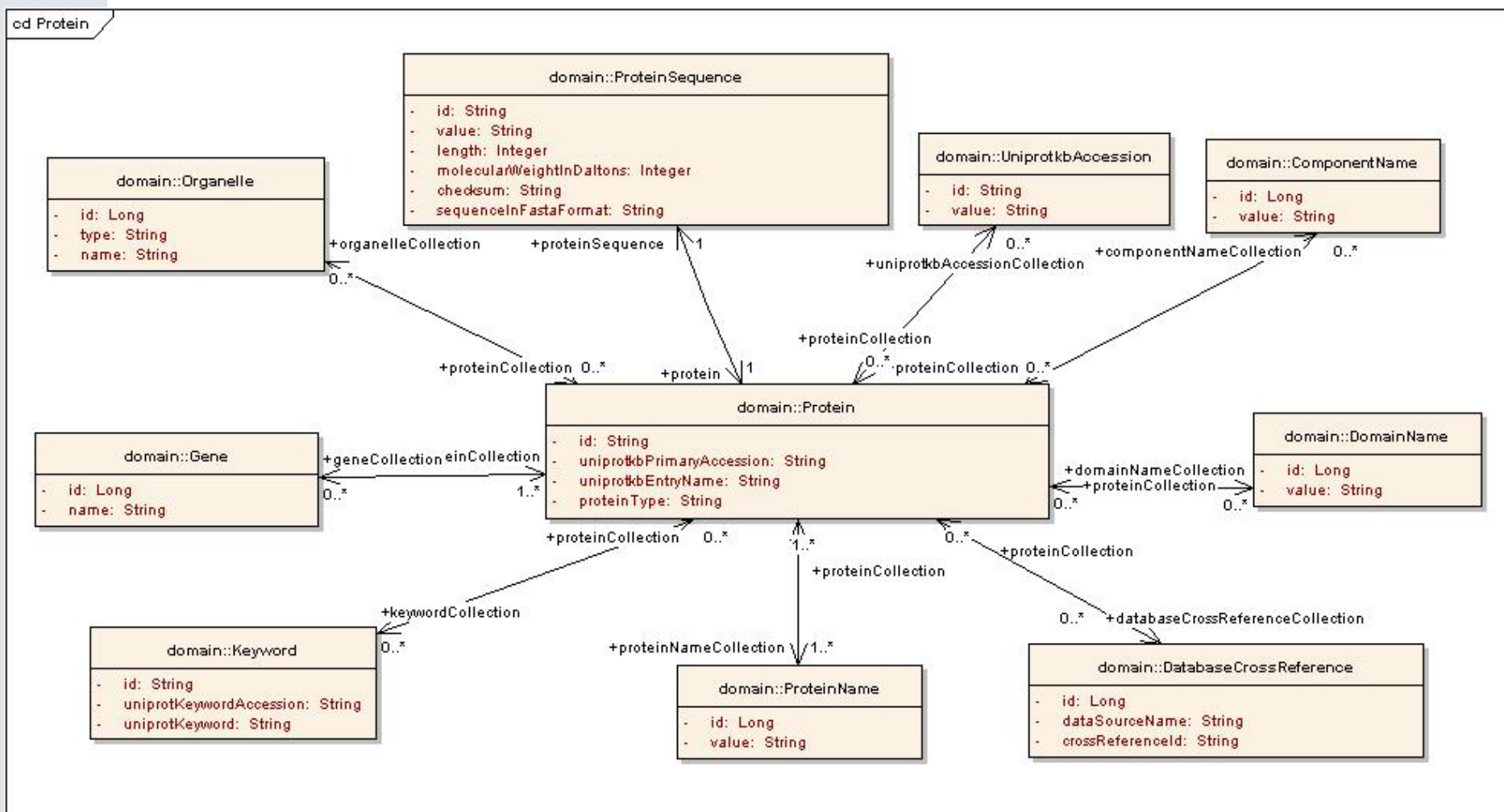
# Data Model

- ▸ In collaboration with Lewis Frey and George Komatsoulis

- ▸ Considerations:

  - • Scientific meaning– Don't do record modeling

  - • Use cases – Consider search criteria objects

  - • caCORE SDK constraints – Consider naming conventions, "id" attribute constraints, supported collection types. e.g. "List" is not supported

  - • Data related constraints – Include only associations or objects based on your data. e.g. Gene to Protein, but not Protein to DNASequence

  - • Semantics – Express semantics and avoid using type attributes. e.g. ProteinFeature subclasses, Lineage
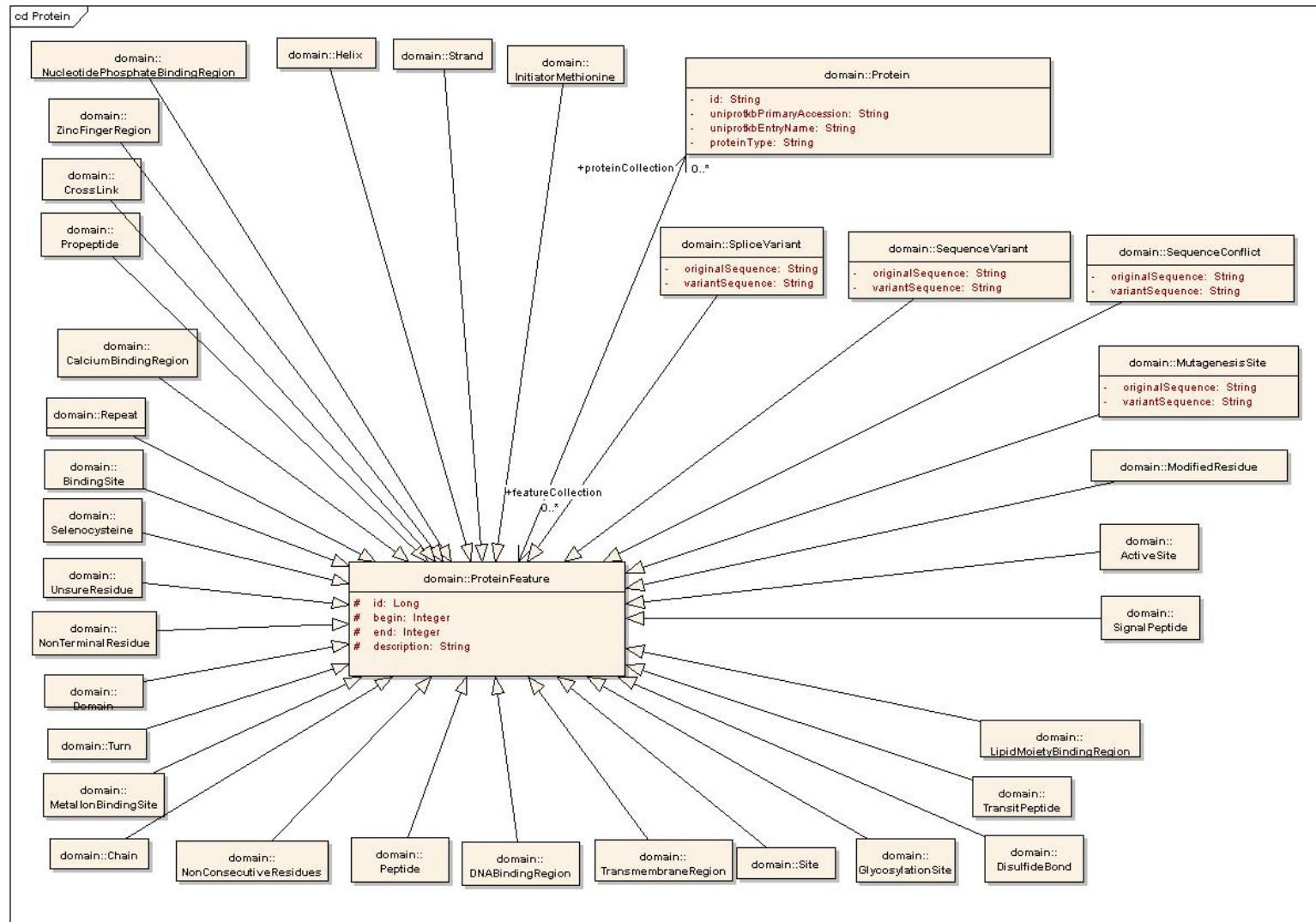
  - • Other projects – Review caCORE/caBIO models

# Data Model
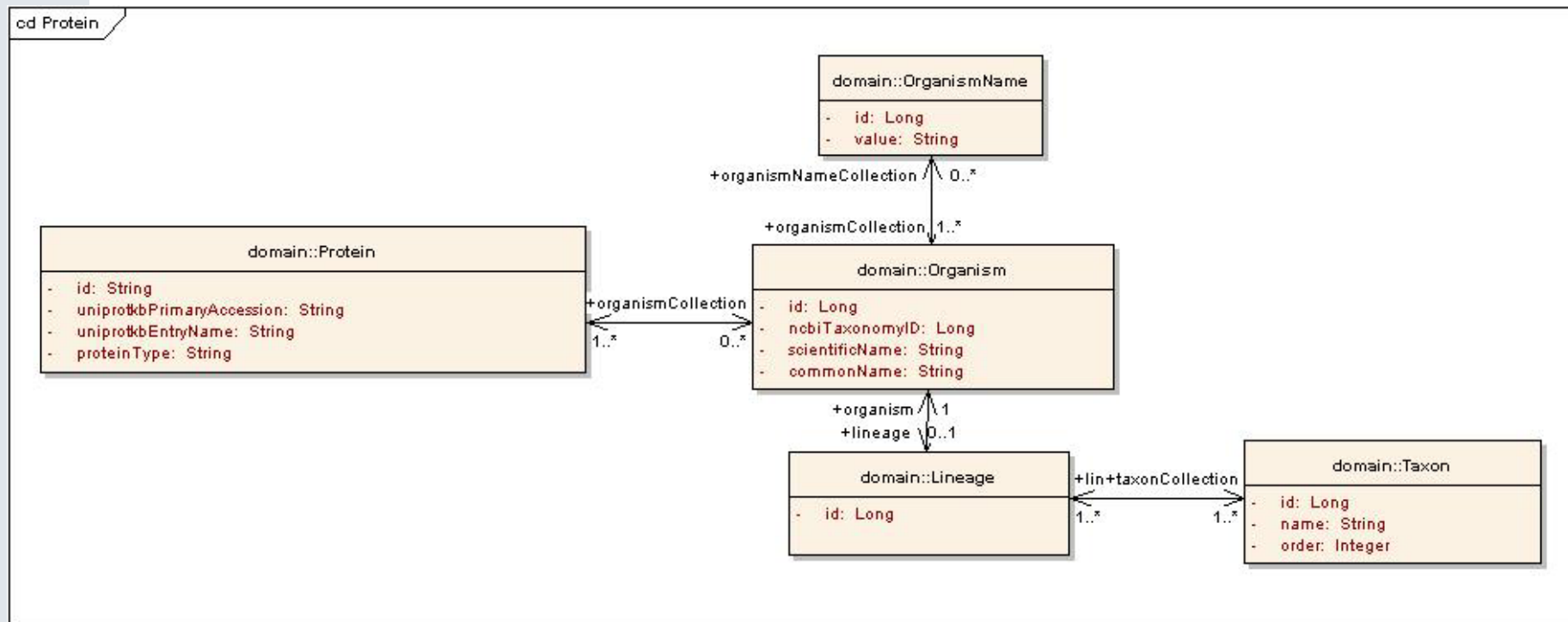
▸ Total 48 objects , 51 attributes

# Data Model

▶ Protein/Gene related objects

# Data Model

▸ Annotation related objects: Protein Features

# Data Model

▸ Taxonomy related objects (Proposed as Taxonomy CDE)

# Semantic Annotation

▸ 149 concepts are used

▸ Loaded to caDSR production server on August 8

▸ Example: Gene.name

–  Property:

  • C42614: Name: The words or language units by which a thing is known.

–  PropertyQualifier1:

  • C43568: Gene_Symbol: A unique gene name approved by an organism specific nomenclature committee.

# Semantic Annotation
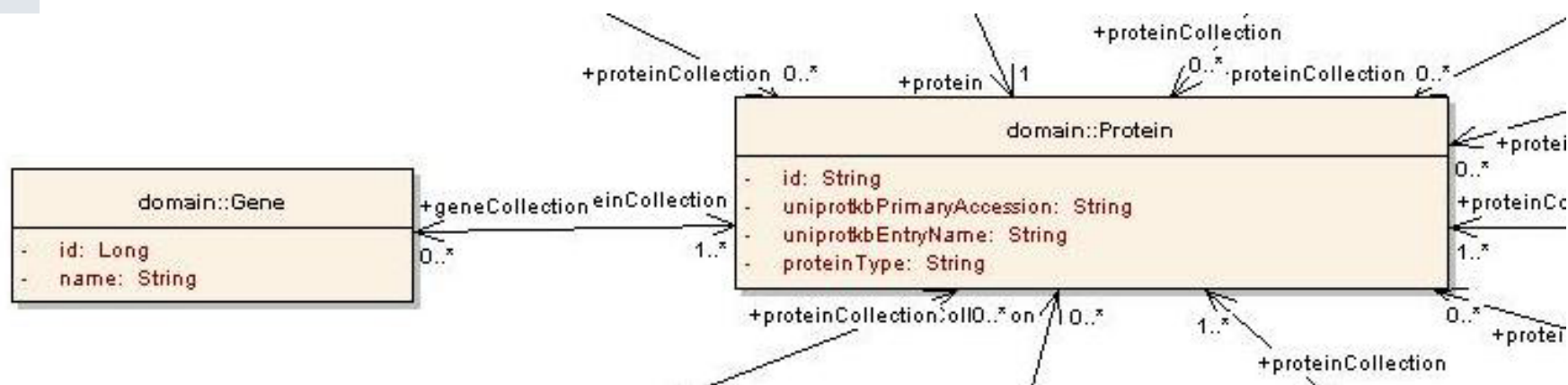
▸ Example: Protein.uniprotkbPrimaryAccession

- Property:

  - C15402: Accession_Number: A control number unique to an object, used to identify it among the other objects in a collection.

- PropertyQualifier1:

  - C4785: UniProt_KB: The UniProt Knowledgebase (UniProtKB), a product of the UniProt consortium, provides a central database of protein sequences with accurate, consistent, rich sequence and functional annotation. The UniProt Knowledgebase consists of two sections: Swiss-Prot - a section containing manually-annotated records with information extracted from literature and curator-evaluated computational analysis, and TrEMBL - a section with computationally analyzed records that await full manual annotation.

- PropertyQualifier2:

  - C25251 : Primary: Occurring first in time or sequence; original; of greatest rank or importance or value.

▶ Retrieve the proteins for gene "BRCA2" (Breast Cancer Gene 2)

```
<caBIGXMLQuery name="testGene2Protein">
 <Target name="edu.georgetown.pir.domain.Protein">
  <Objects name="edu.georgetown.pir.domain.Gene">
   <Property name="name" predicate="equal" value="BRCA2"/>
  </Objects>
 </Target>
</caBIGXMLQuery>
```
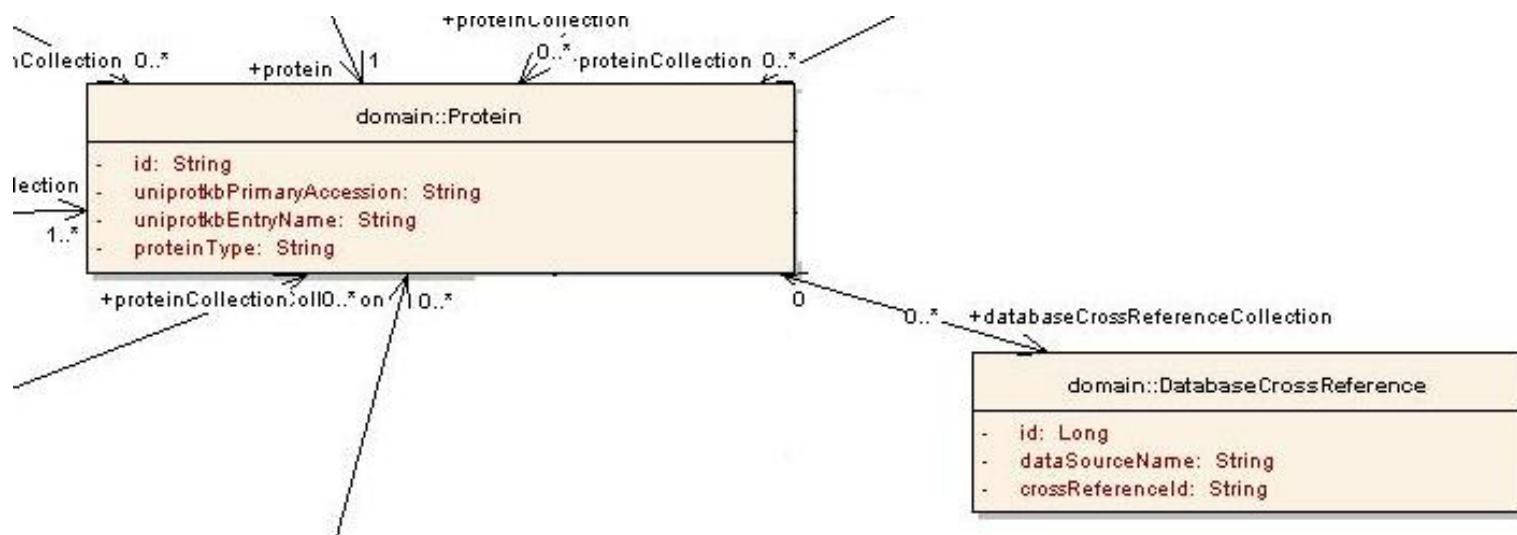
# Using PIR Grid Service

▶ Find all the proteins that contain the domain "BRCA2 repeat" (PFAM:PF00634, a domain in Breast cancer type 2 susceptibility protein)
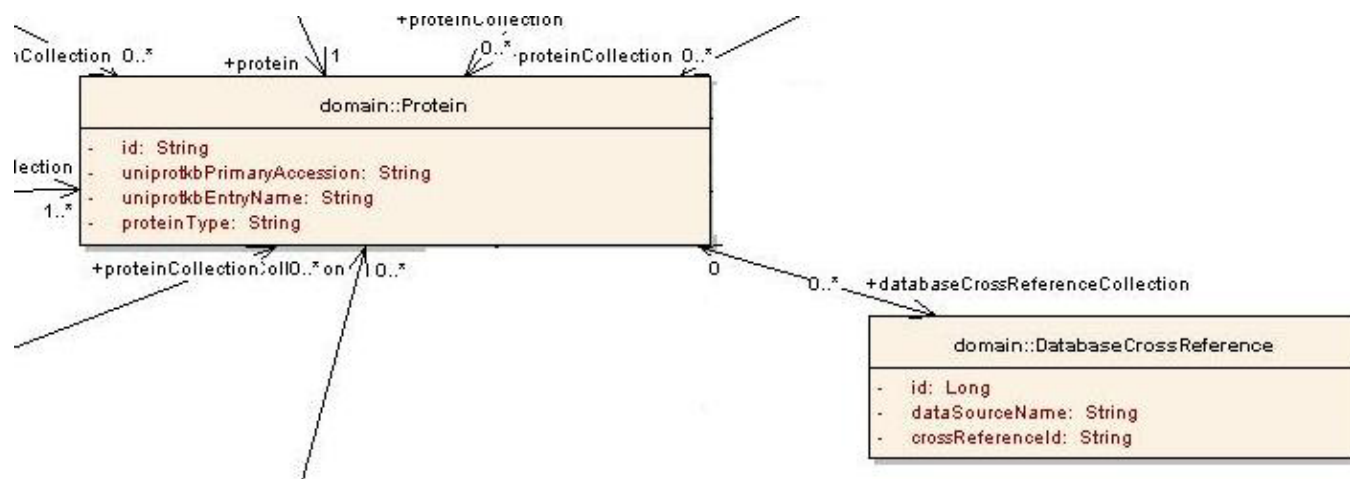
```xml
<caBIGXMLQuery name="testPfam2Protein">
  <Target name="edu.georgetown.pir.domain.Protein">
     <Objects name="edu.georgetown.pir.domain.DatabaseCrossReference">
       <Property name="crossReferenceId" predicate="equal" value="PF00634"/>
     </Objects>
  </Target>
</caBIGXMLQuery>
```

# Using PIR Grid Service

▶ ID mapping: Find all the database cross-references from various databases corresponding to RefSeq Accession NP_061820

```
<caBIGXMLQuery name="testIDMapping">
 <Target name="edu.georgetown.pir.domain.DatabaseCrossReference"
    path="edu.georgetown.pir.domain.Protein">
  <Objects name="edu.georgetown.pir.domain.DatabaseCrossReference">
   <Property name="dataSourceName" predicate="equal"  value="RefSeq"/>
   <Property name="crossReferenceId" predicate="equal" value="NP_061820"/>
  </Objects>
</Target>
</caBIGXMLQuery>
```
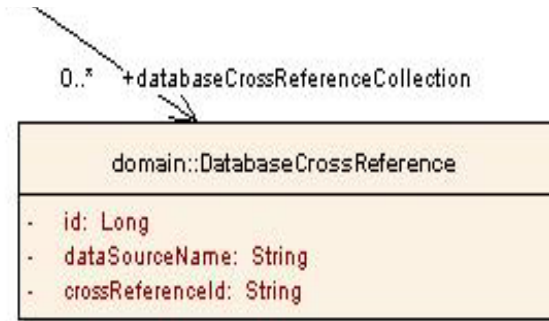
▸ caBIO and PIR databaseCrossReference objects and dataSourceName – Flexibility vs. better semantics

**Current**



**Future (?)**